

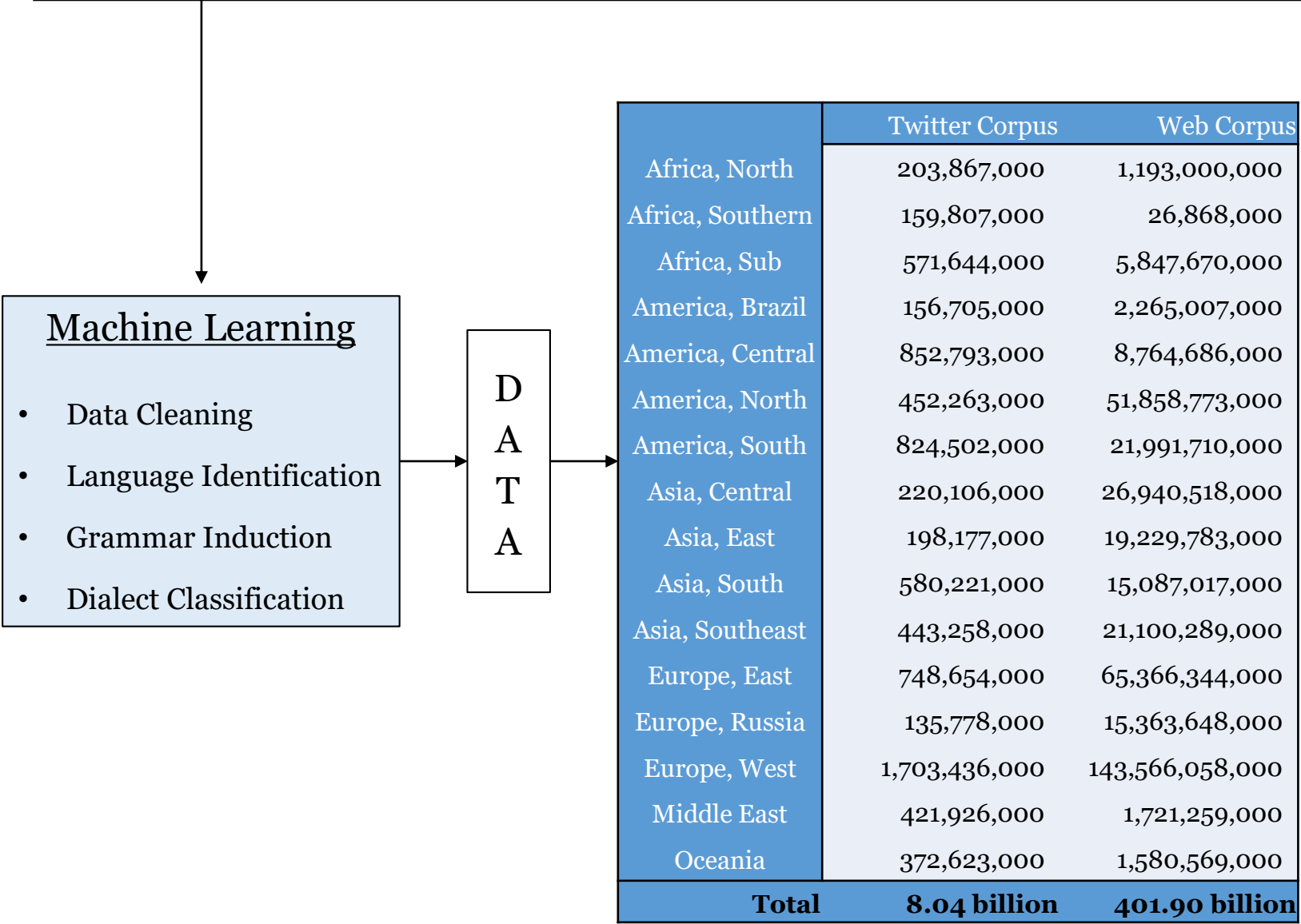
Traditional approaches require *asking* people once every decade what languages they use.

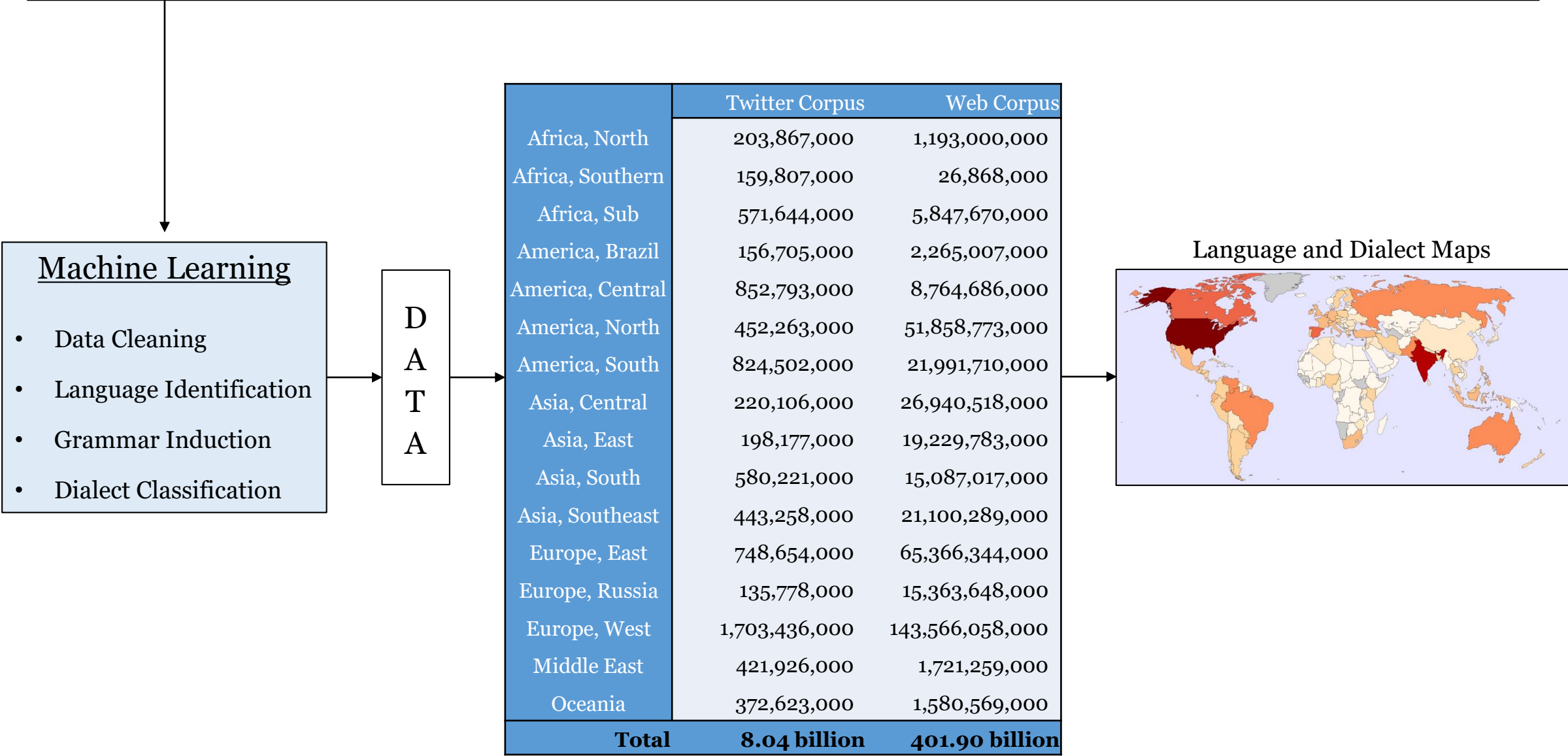
What if we could **speed up** and **scale up** this survey process?



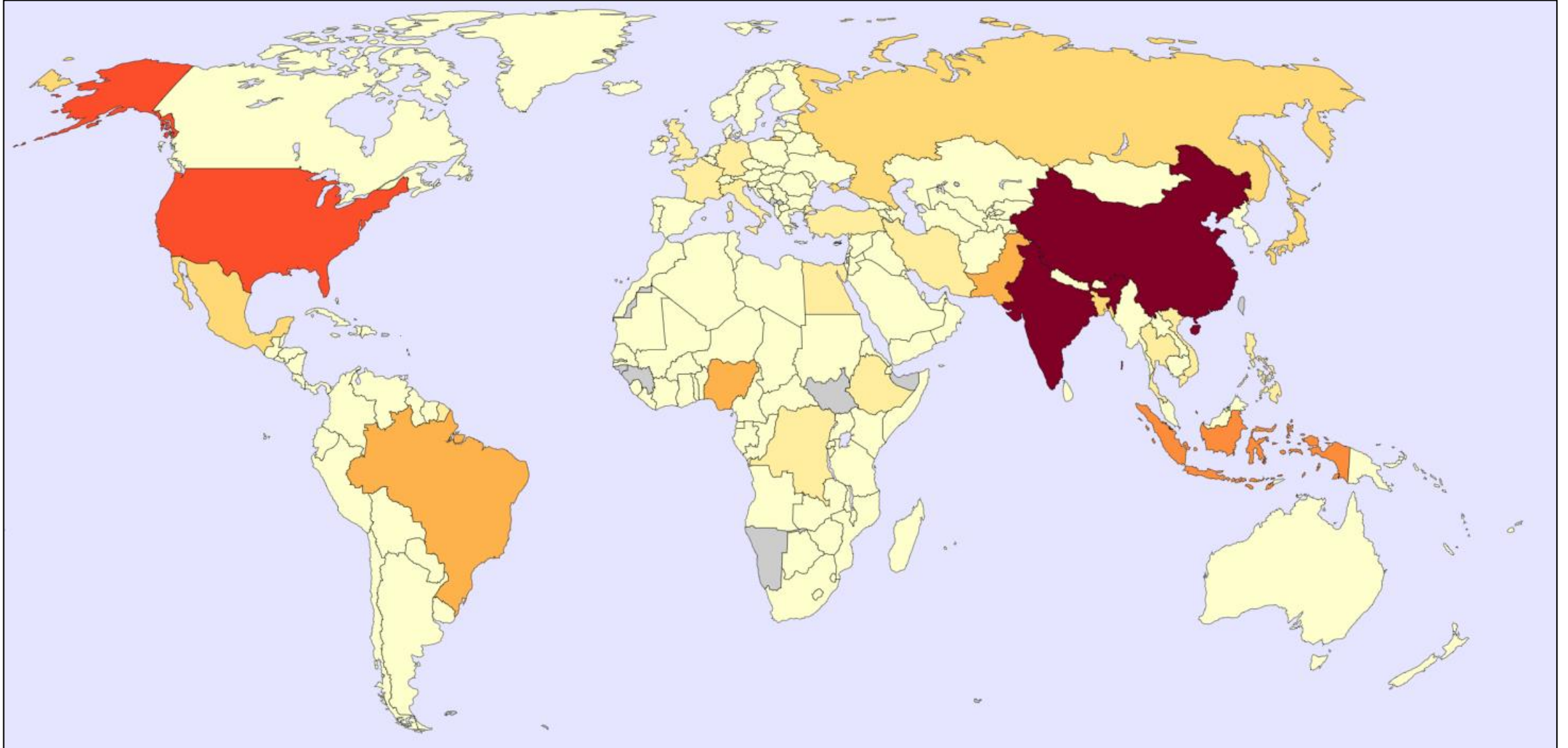
### Machine Learning

- Data Cleaning
- Language Identification
- Grammar Induction
- Dialect Classification

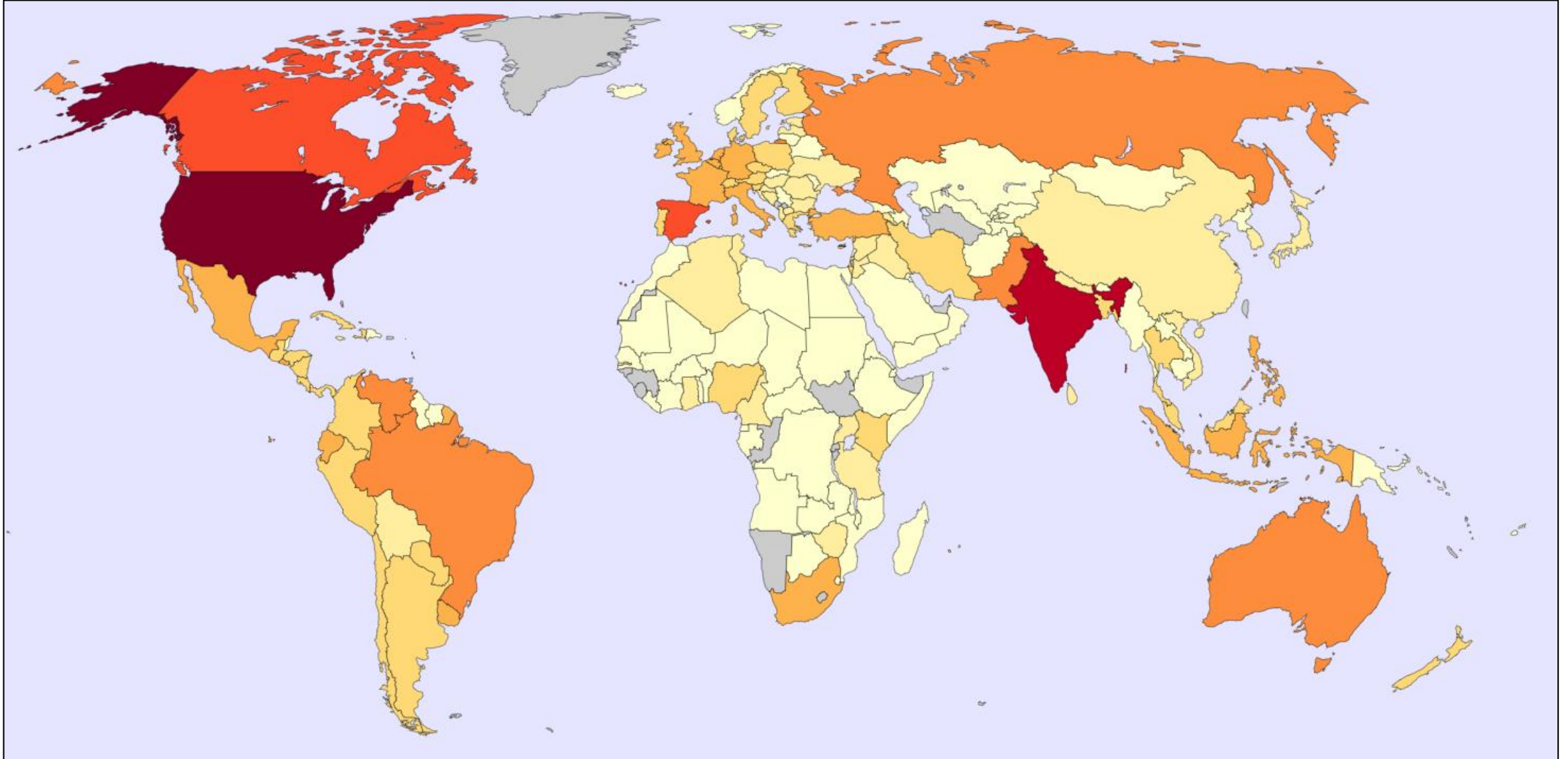




Population: Where do people live?

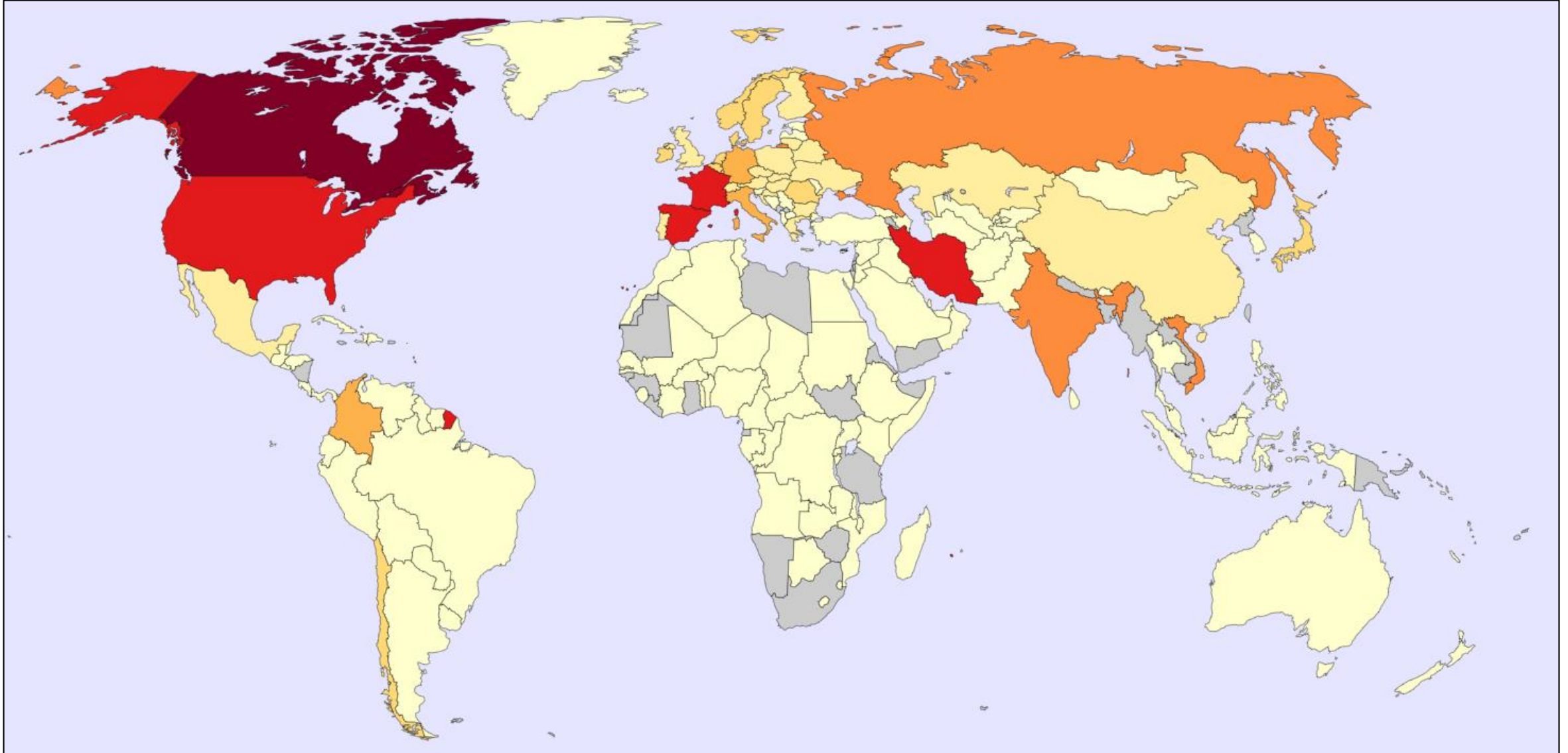


Twitter Density: Where do people Tweet?

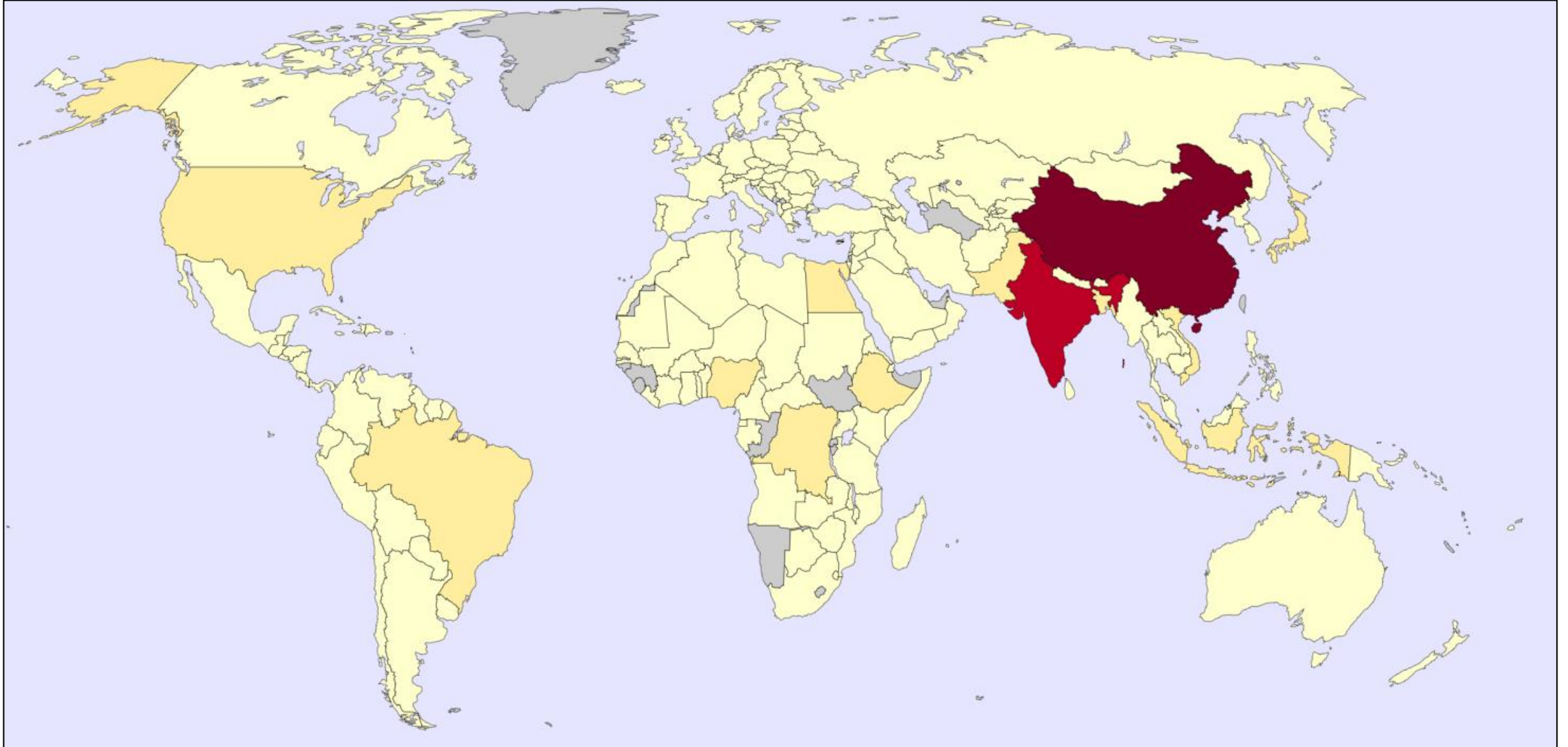




Web Density: Where do web documents come from?

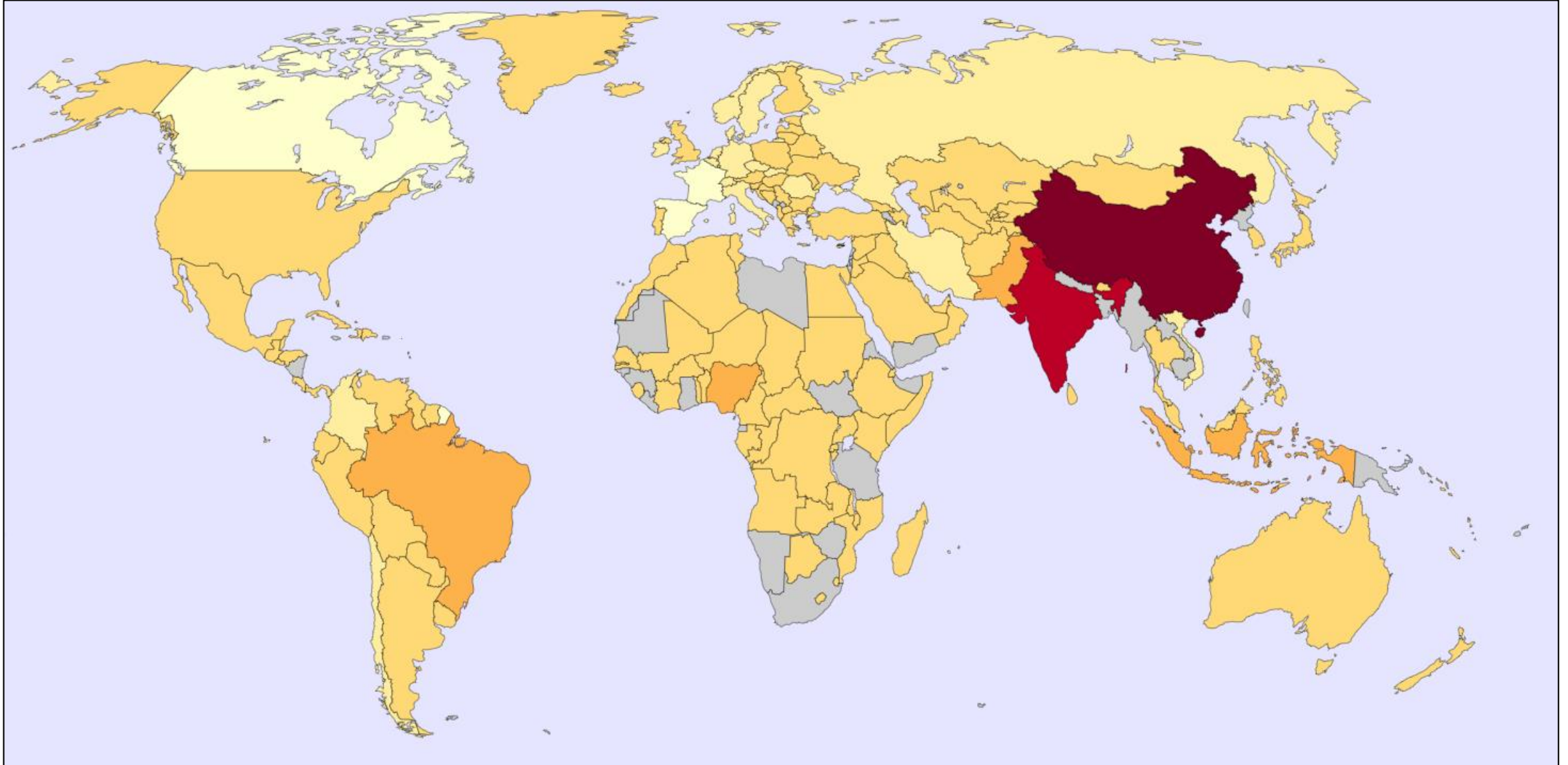


Where is Twitter not used? (i.e., under-represented countries)

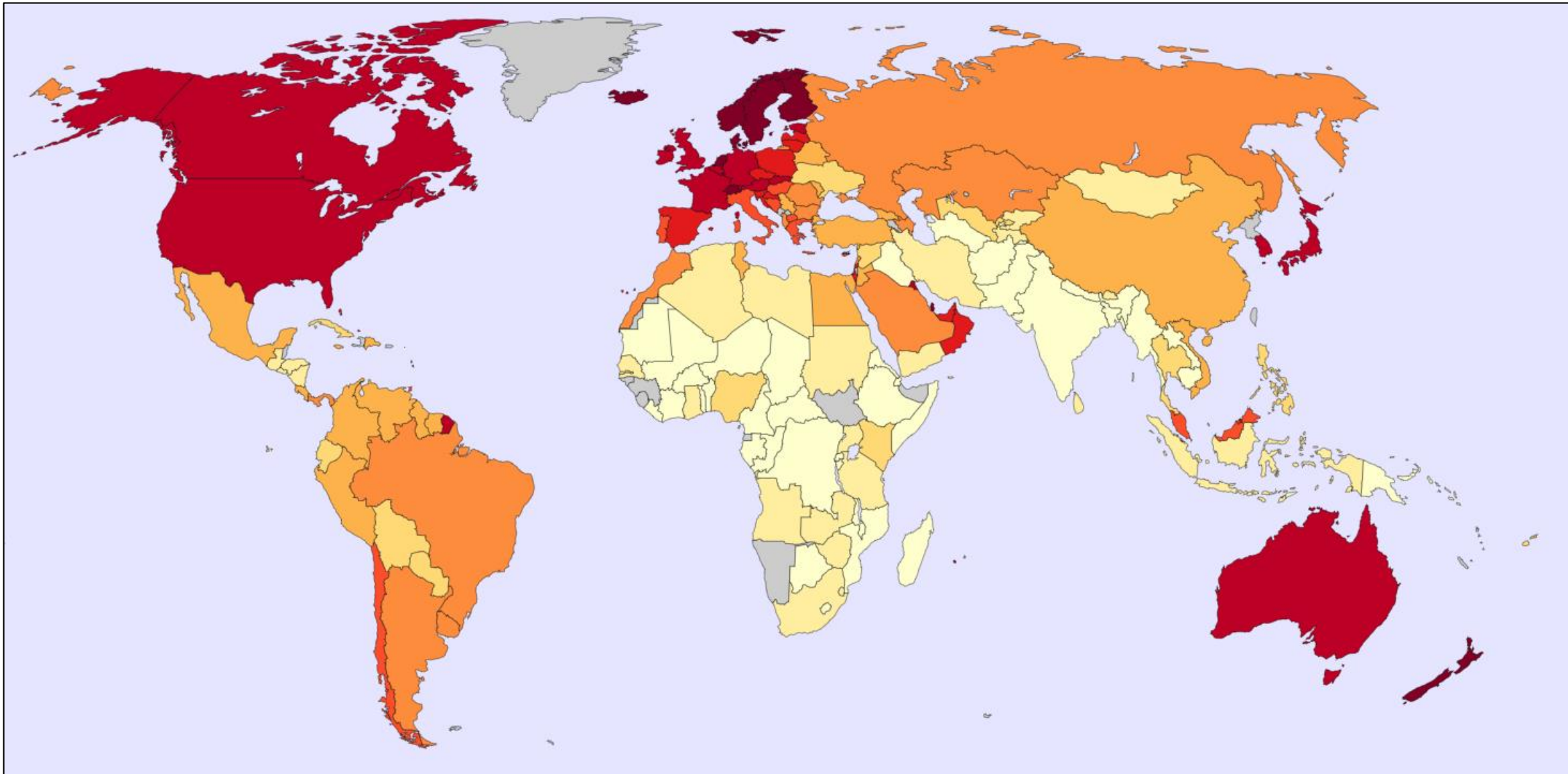




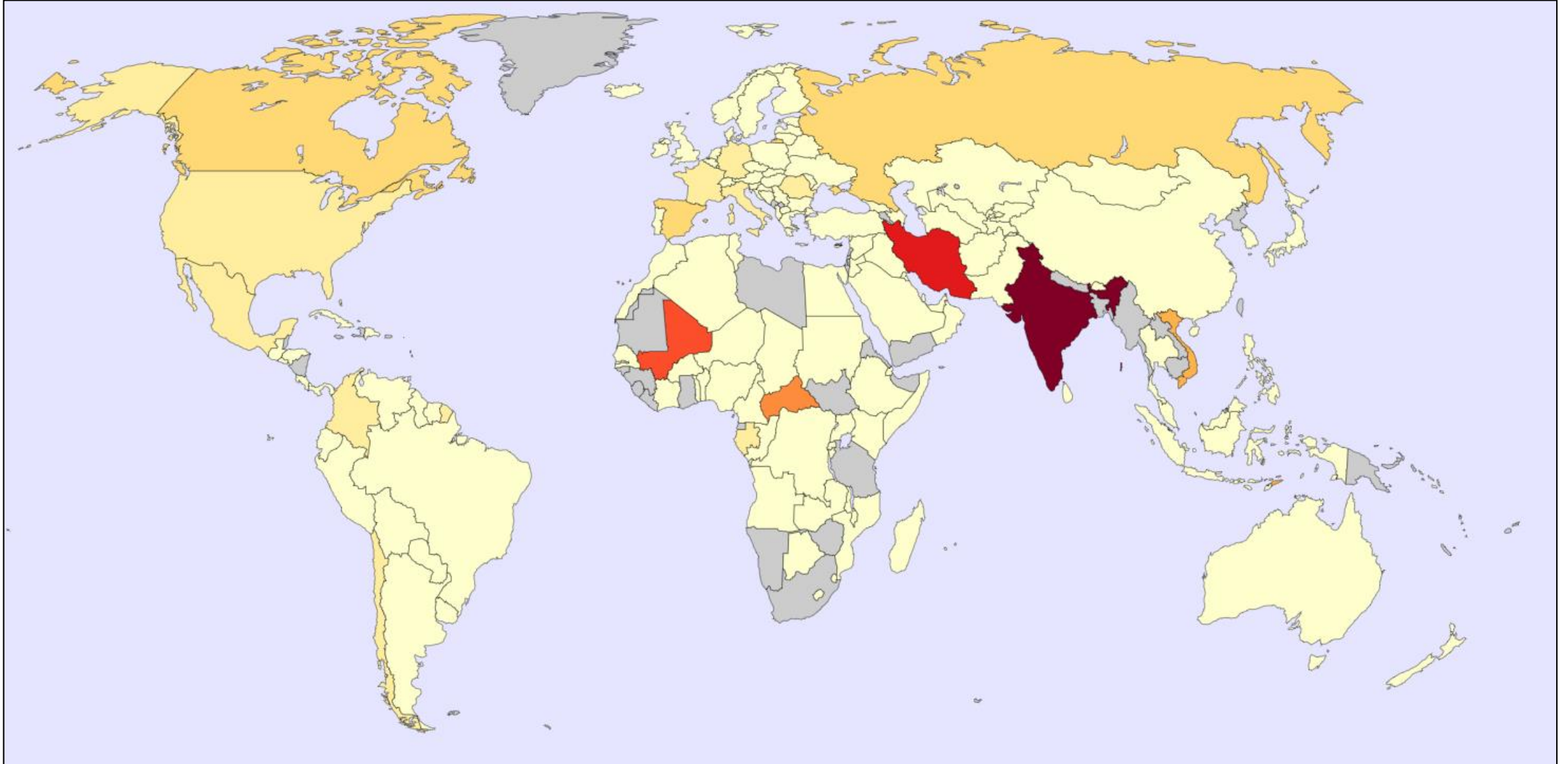
Where are web documents scarce? (i.e., under-represented countries)



Internet Usage: Does the percentage of the population with internet access influence these datasets?

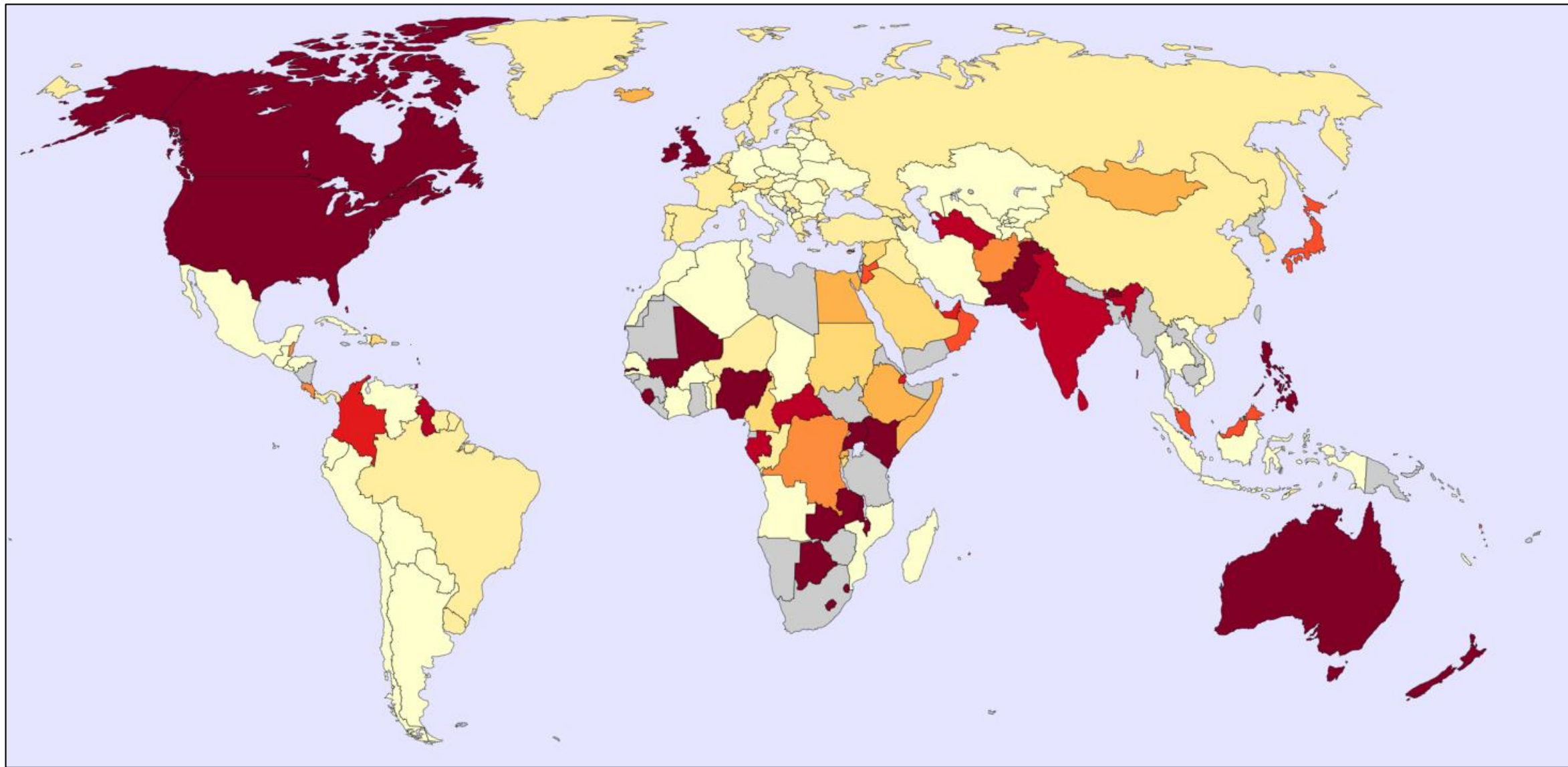


## What would the web corpus look like if everyone had internet access?

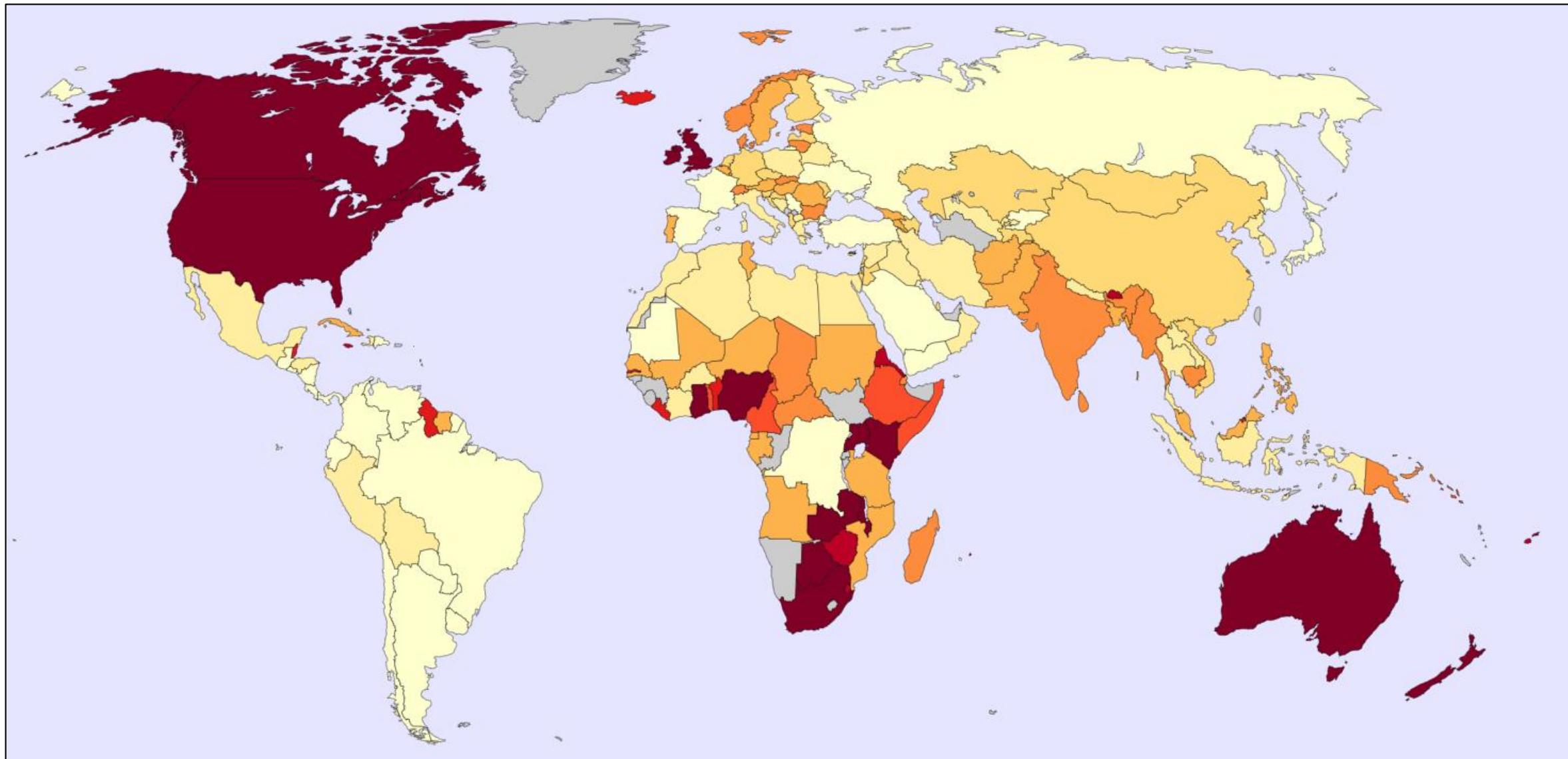




Countries by their percent usage of English (web): Darker red means more monolingual usage

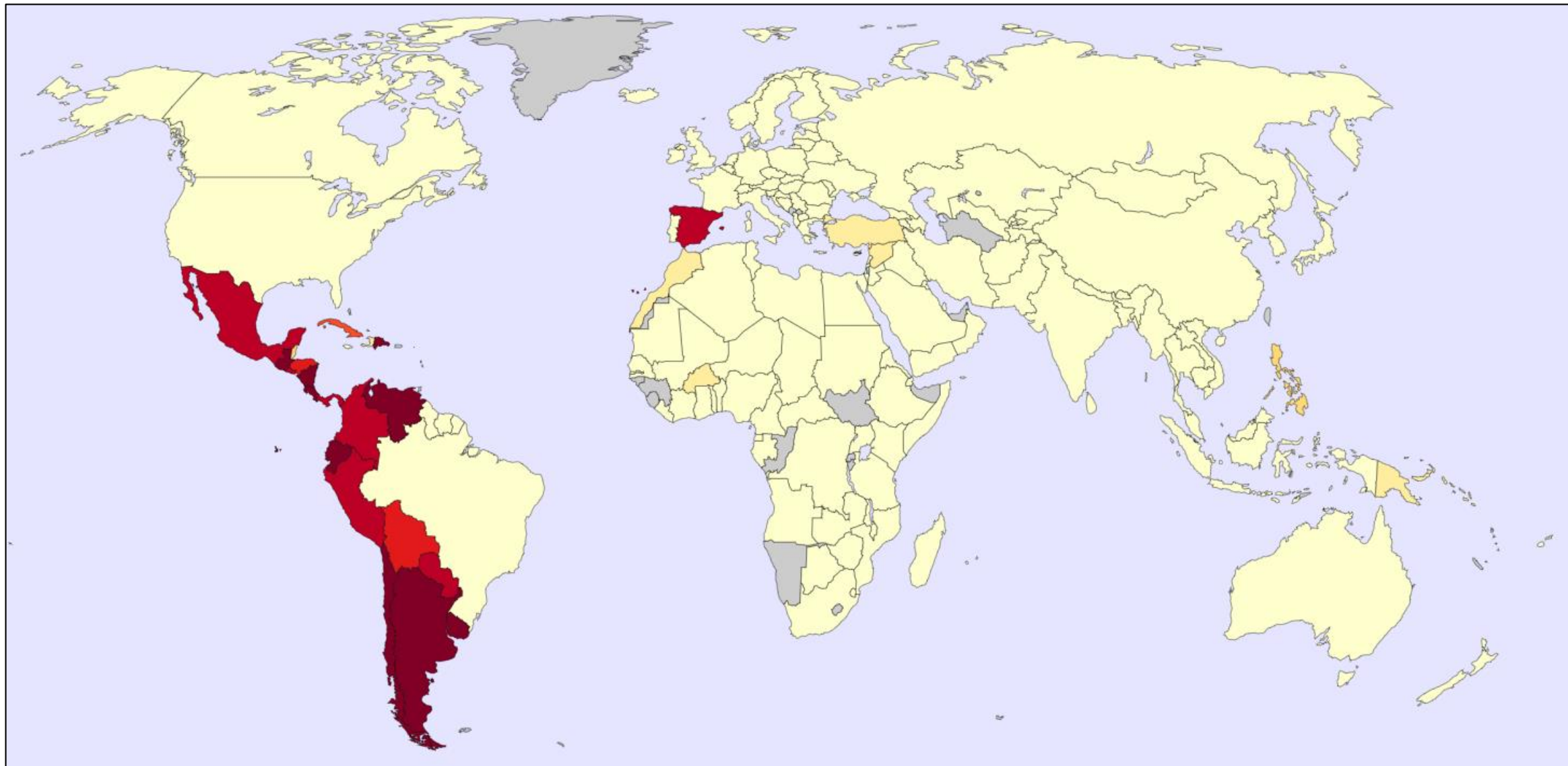


Countries by their percent usage of English (Twitter): Do the datasets capture different populations?

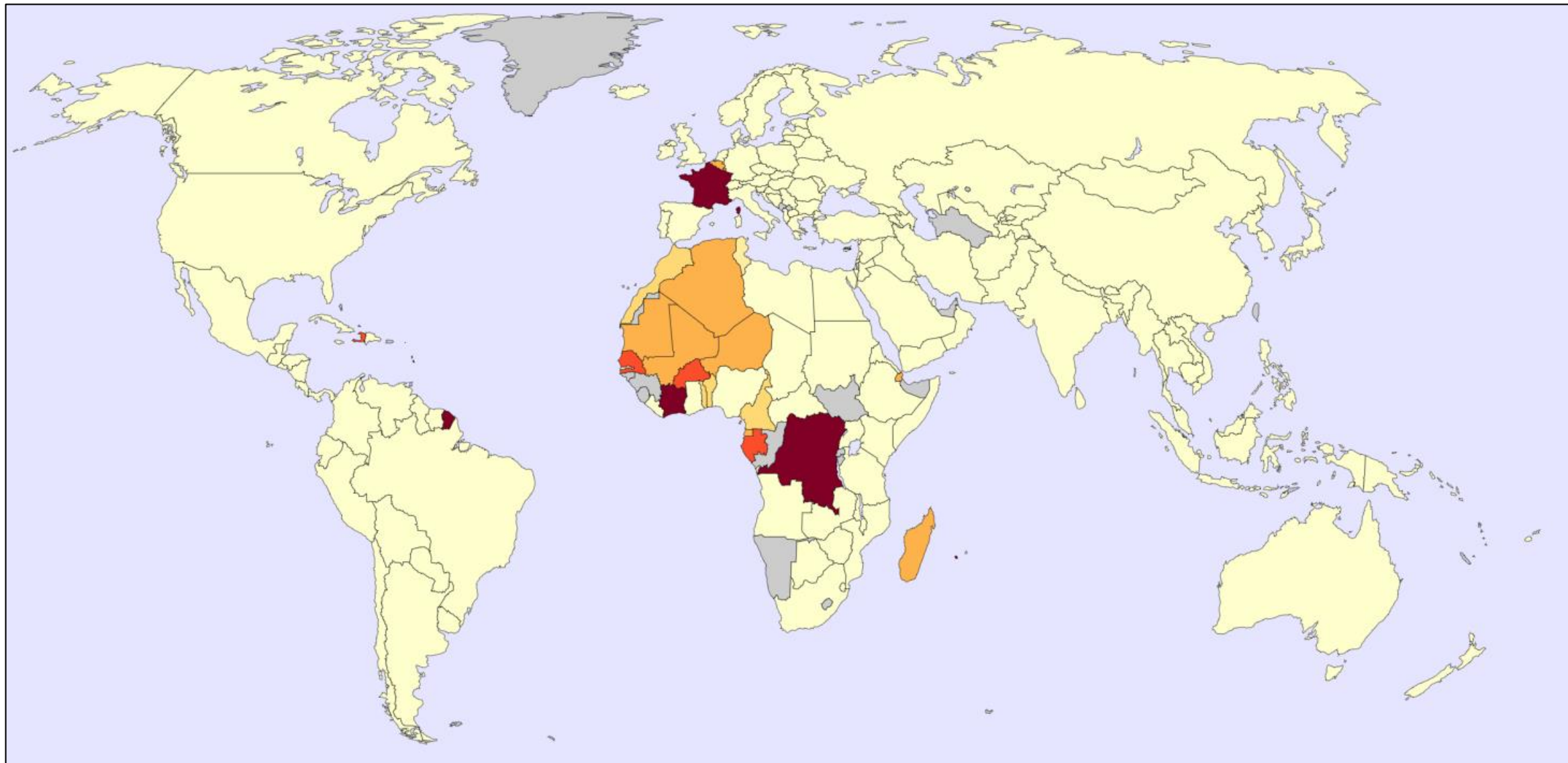




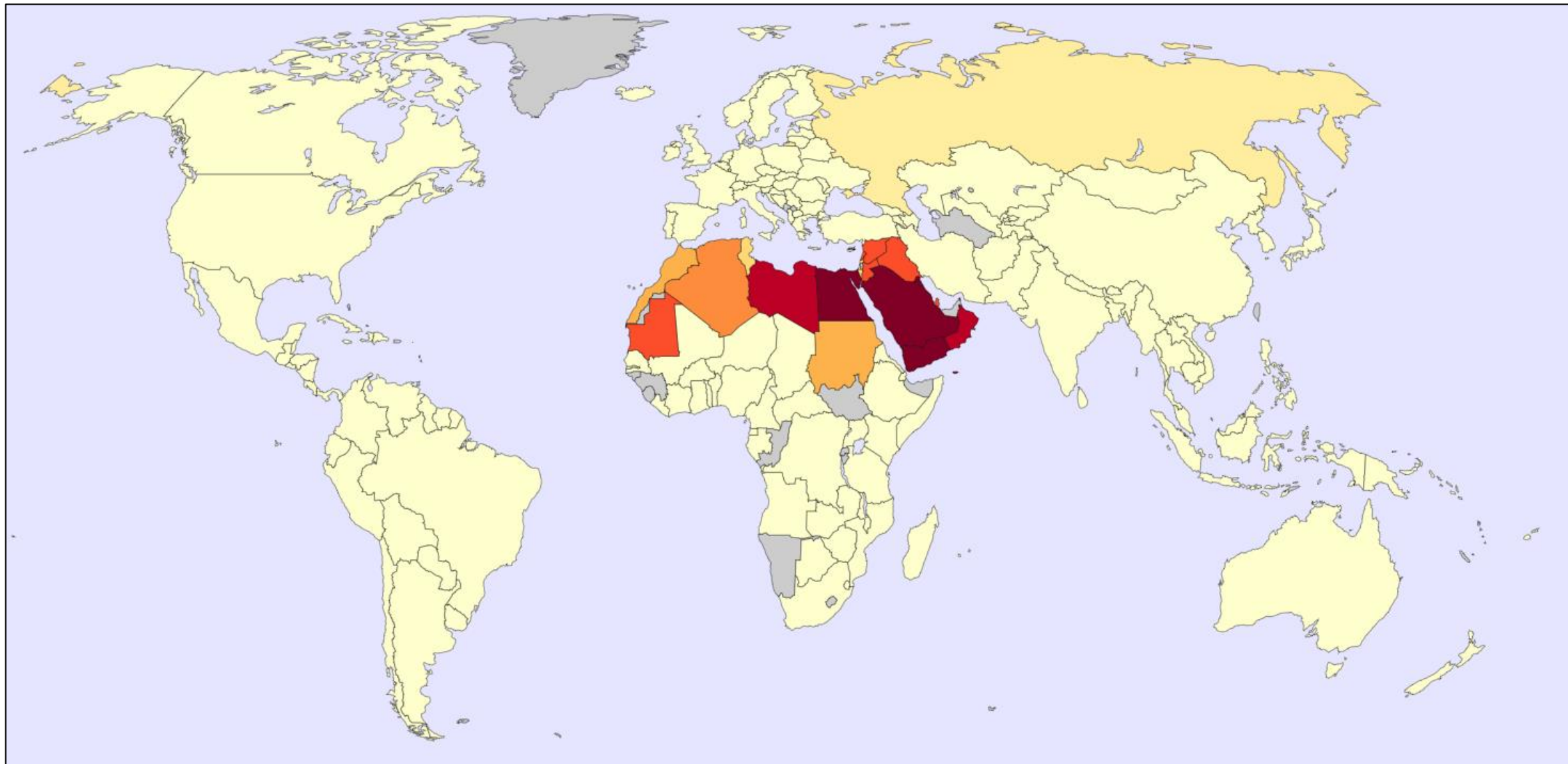
Countries by their percent usage of Spanish (Twitter): English is a global language, but not Spanish



Countries by their percent usage of French (Twitter): French is also widely used

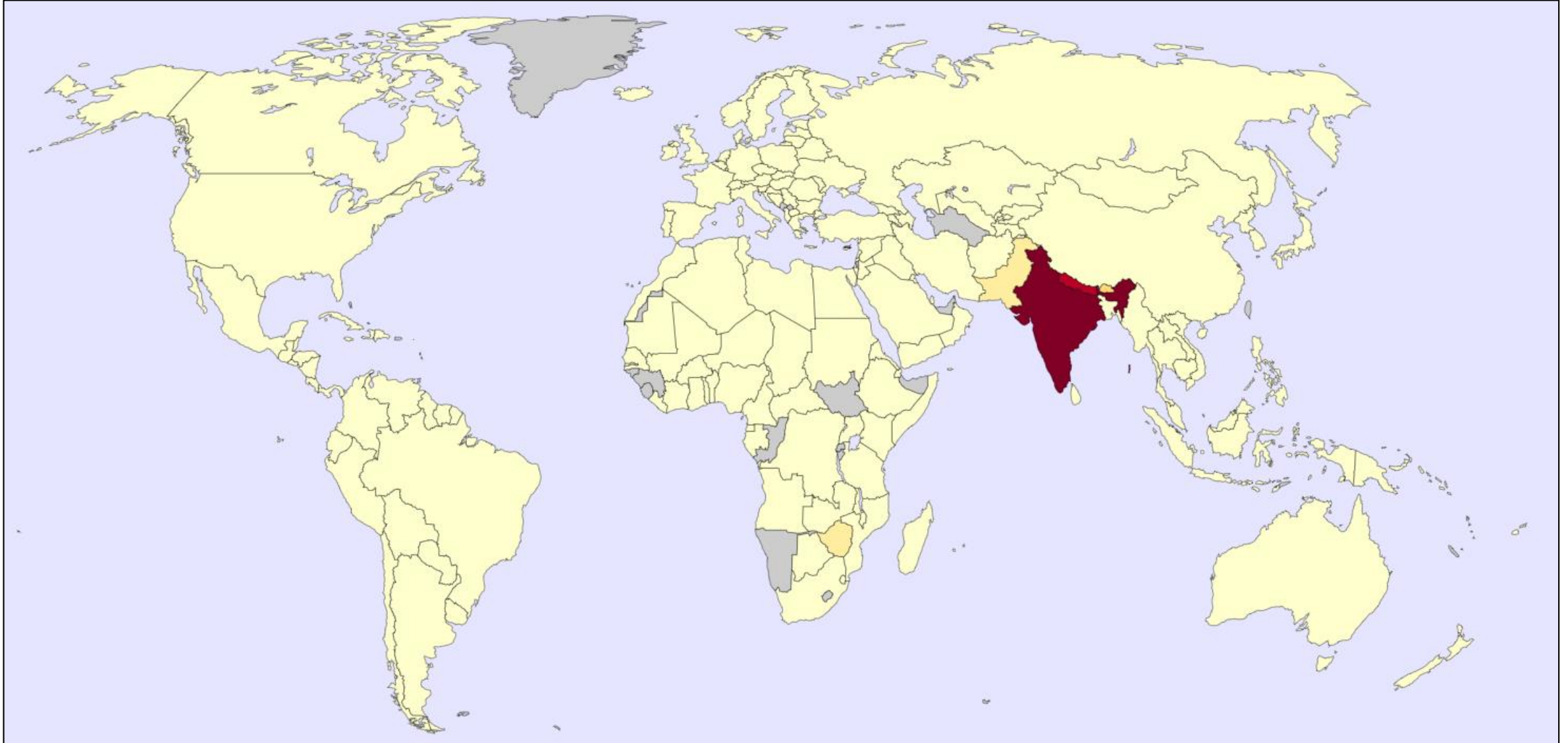


Countries by their percent usage of Arabic (Twitter): Some languages move with along moving populations





## Countries by their percent usage of Hindi (Twitter): But others are restricted geographically

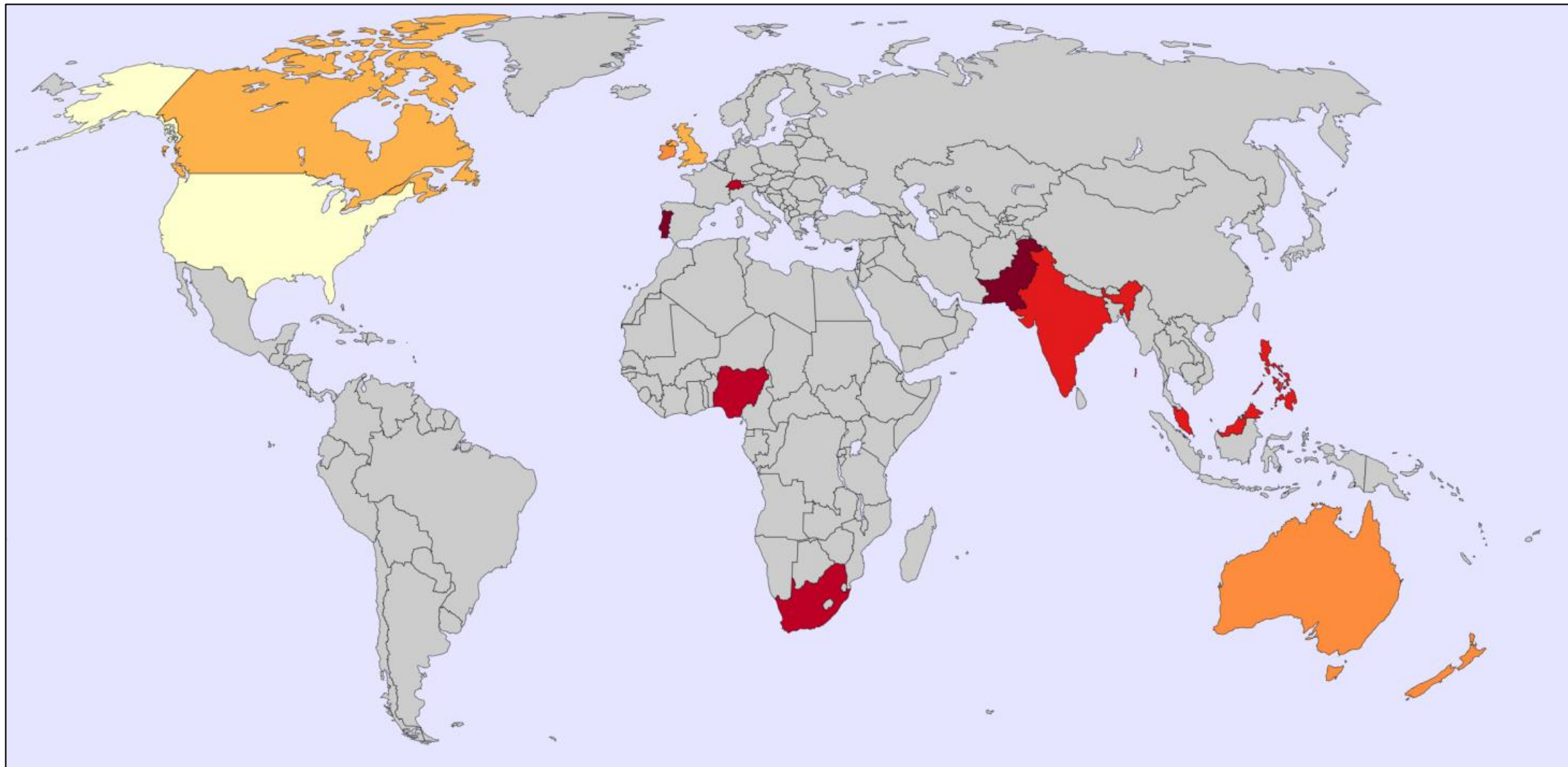


Countries by their percent usage of Thai (Twitter): Here, Thai is well-represented... but only in a narrow region





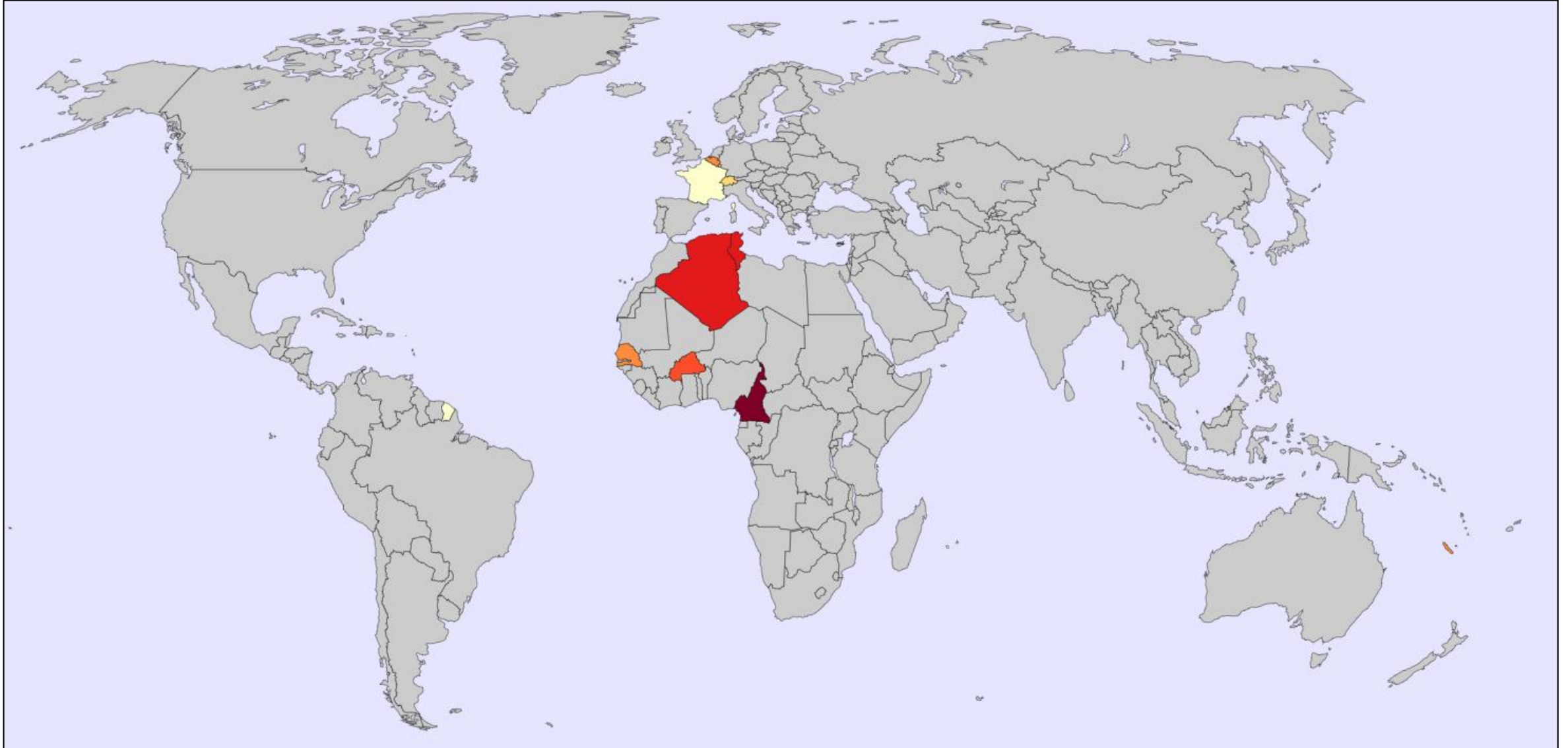
Uniqueness of English Dialects (Web): We can go beyond surveys by modelling the data from each country



Uniqueness of Spanish Dialects (Web): The dialect model shows the Spain is the central variety of Spanish

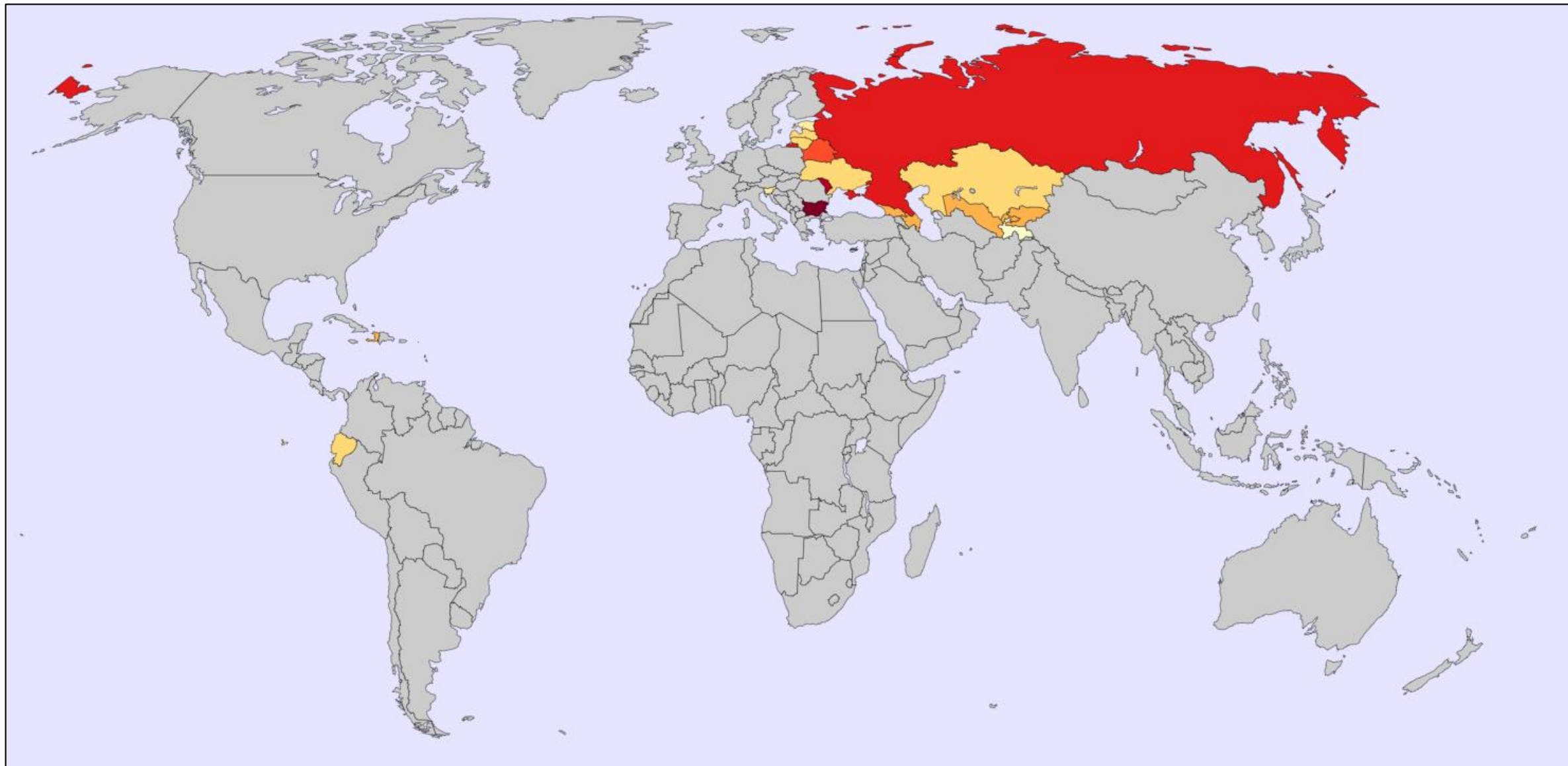


Uniqueness of French Dialects (Web): There aren't that many country-level dialects of French





Uniqueness of Russian Dialects (Web): Unlike other major languages, Russian is restricted to a (large) contiguous region



Explore the data at [earthLings.io](https://earthlings.io)